



1 Introduction

This Brook+ sample computes the Cholesky factorization of a positive definite matrix. A positive definite matrix, \mathbf{A} , can be represented as a product of a lower triangular matrix and its transpose:

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T$$

For a 2x2 matrix,

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 \\ l_{21} & l_{22} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} \\ 0 & l_{22} \end{bmatrix}$$

Simplifying the RHS we get,

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} l_{11}^2 & \\ l_{21} * l_{11} & \end{bmatrix} \begin{bmatrix} l_{11} * l_{21} \\ l_{21} * l_{11} + l_{22}^2 \end{bmatrix}$$

The entries of \mathbf{L} can be computed as follows:

1. $l_{11} = \text{sqrt}(a_{11})$
2. $l_{21} = a_{21} / l_{11}$
3. $l_{22} = \text{sqrt}(a_{22} - l_{21}^2)$

The last step can be separated into:

3.a) $a_{22} = (a_{22} - l_{21}^2)$

3.b) $l_{22} = \text{sqrt}(a_{22})$

These steps for a 2 x 2 matrix can be generalized for any positive-definite matrix of dimensions $m \times n$ as an iterative algorithm consisting of three basic steps:

1. square-root,
2. normalization, and
3. sub-matrix update.

These are applied iteratively to progressively smaller regions of the original matrix, \mathbf{A} . In the naïve version, the algorithm processes individual elements of the matrix \mathbf{A} . This algorithm, according to Jung et al. [1], is given below and shown in Figure 1.

```

for k = 1 to n-1 do
  A(k, k) = sqrt(A(k, k)) // Square root Step
  A(k+1:n, k) = A(k+1:n, k)/A(k, k) // Normalization Step
  for j = k+1 to n do
    A(j, k+1) = A(j, k+1) - A(j, 1:k)T x A(k+1, 1:k) // Sub-matrix update step
  end for
end for
A(n, n) = sqrt(A(n, n))

```

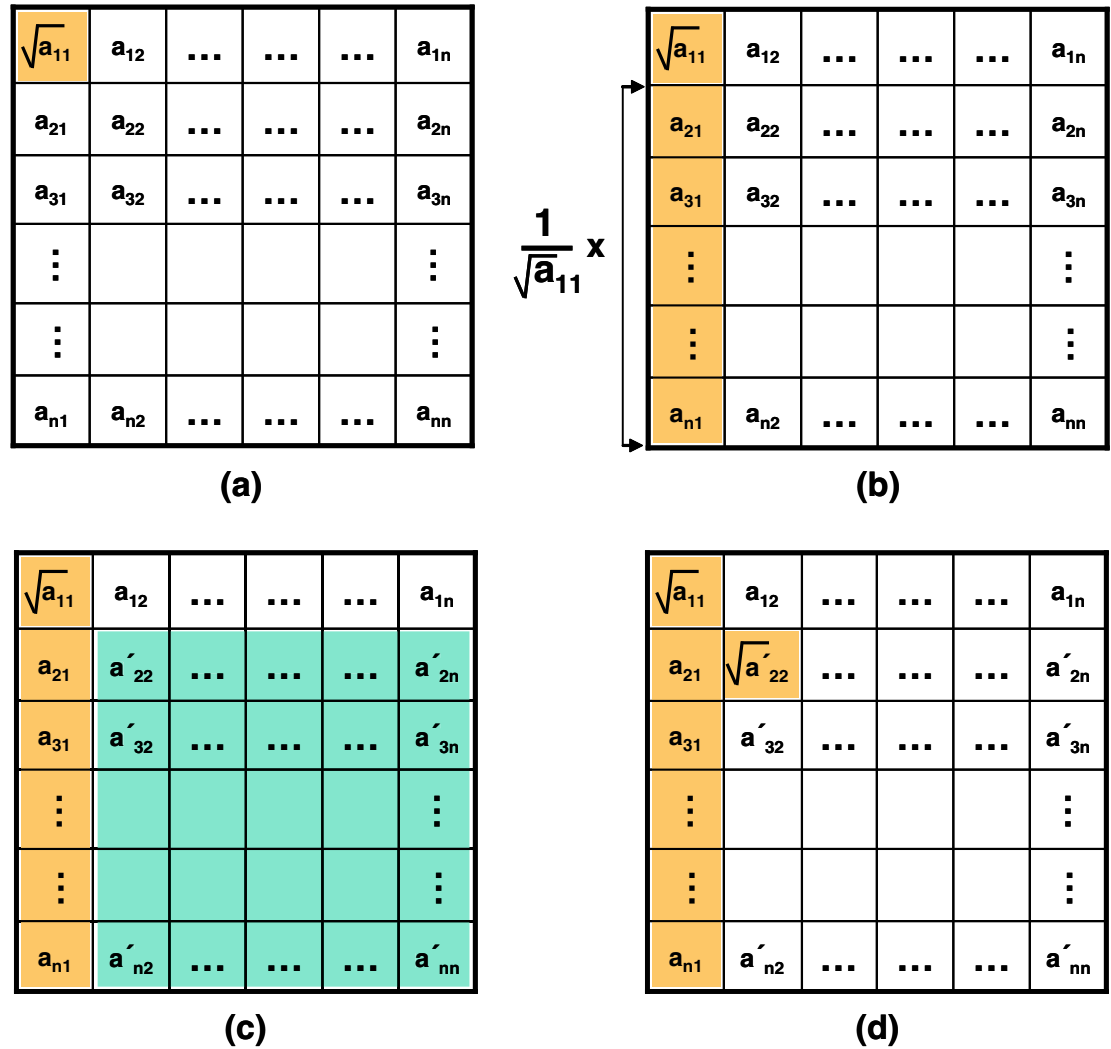


Figure 1 Elementwise Computation of the Cholesky Factorization

For Figure 1:

- Square root, $l_{11} = \sqrt{a_{11}}$.
- Normalize a_{21} (highlighted in black).
- Update sub-matrix a_{22} (highlighted in black).
- Repeat steps on sub-matrix starting with square-root.

2 Blocked Cholesky Factorization

The elementwise algorithm is inefficient because it processes only one column in every iteration. We can achieve a significant performance gain by considering entire matrix blocks at a time. In block-partitioned form, the Cholesky factorization of a positive matrix \mathbf{A} of size $m \times m$ is:

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12}^T \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{11} & \mathbf{0} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{11}^T & \mathbf{L}_{21}^T \\ \mathbf{0} & \mathbf{L}_{22}^T \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{L}_{11}\mathbf{L}_{11}^T & \mathbf{L}_{11}\mathbf{L}_{21}^T \\ \mathbf{L}_{21}\mathbf{L}_{11}^T & \mathbf{L}_{21}\mathbf{L}_{21}^T + \mathbf{L}_{22}\mathbf{L}_{22}^T \end{bmatrix}$$

This shows that \mathbf{L}_{11} is the Cholesky factor of \mathbf{A}_{11} . Assume that a routine Cholesky computes the Cholesky factorization of a matrix using the elementwise computation described earlier. We then can have the following iterative algorithm for block-partitioned Cholesky factorization of a matrix:

1. Square Root Step: $\mathbf{L}_{11} = \text{Cholesky}(\mathbf{A}_{11})$
2. Normalization Step: $\mathbf{L}_{21} = \mathbf{A}_{21} (\mathbf{L}_{11}^T)^{-1}$
3. Sub-matrix Update Step: $\mathbf{A}_{22}' = \mathbf{A}_{22} - \mathbf{L}_{21} \mathbf{L}_{21}^T = \mathbf{L}_{22} \mathbf{L}_{22}^T$
4. Repeat steps 1 to 3 on residual matrix \mathbf{A}_{22}' .

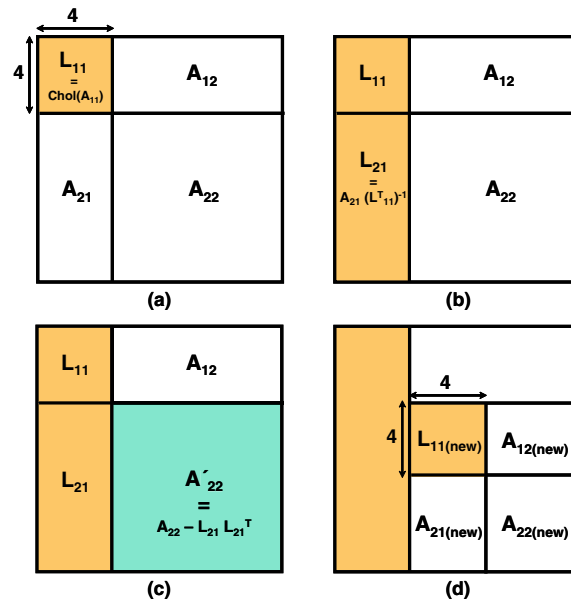


Figure 2 Blocked Cholesky Factorization

For [Figure 2](#):

- a. Block partitions of matrix \mathbf{A}
- b. Square Root Step: $\mathbf{L}_{11} = \text{Cholesky}(\mathbf{A}_{11})$
- c. Normalization Step: $\mathbf{L}_{21} = \mathbf{A}_{21} (\mathbf{L}_{11}^T)^{-1}$
- d. Sub-matrix Update Step: $\mathbf{A}_{22}' = \mathbf{A}_{22} - \mathbf{L}_{21} \mathbf{L}_{21}^T = \mathbf{L}_{22} \mathbf{L}_{22}^T$

For the second pass, the input to the algorithm is the updated sub-matrix \mathbf{A}_{22}' (yellow region).

3 Brook+ SIMD Implementation

For the Brook+ sample in the SDK, we choose a block size of 4x4. This can be stored conveniently in four float4 variables. We provide a kernel for each step of the algorithm. The square root step is achieved by computing an elementwise Cholesky factorization of the 4x4 matrix. The normalization kernel performs two tasks:

1. It computes the inverse-transpose of the 4x4 Cholesky factor L_{11} .
2. It uses the result to perform the normalization by matrix multiplication.

The sub-matrix update kernel consists of a matrix-multiplication step followed by matrix subtraction.

4 Data Transfer

Our blocked Cholesky factorization sample follows a multi-pass approach that iterates over progressively smaller regions of the input matrix with each pass. Specifically, at the end of each pass, we are only interested in the residual updated sub-matrix A_{22}' (see Figure 2). We use a ping-pong strategy to facilitate data transfer from one pass to the next. Two streams, A and L , are declared; the input matrix is initially copied to both. Brook+ enforces the rule that the same stream cannot be bound simultaneously for reading and writing. We always bind L for writing; thus, the partial results always remains in the L stream. At the end of each pass, we copy the sub-domain in L corresponding to the updated sub-matrix A_{22}' back to A .

5 Performance

The performance of Brook+ Cholesky Factorization was tested on a system with the AMD Athlon™ 64 X2 Dual Core Processor CPU and an ATI Radeon™ HD 4800 Series GPU. Figure 3 shows a comparison of the time taken to compute the Cholesky Factorization of matrices of varying sizes for the GPU versus that of the CPU.

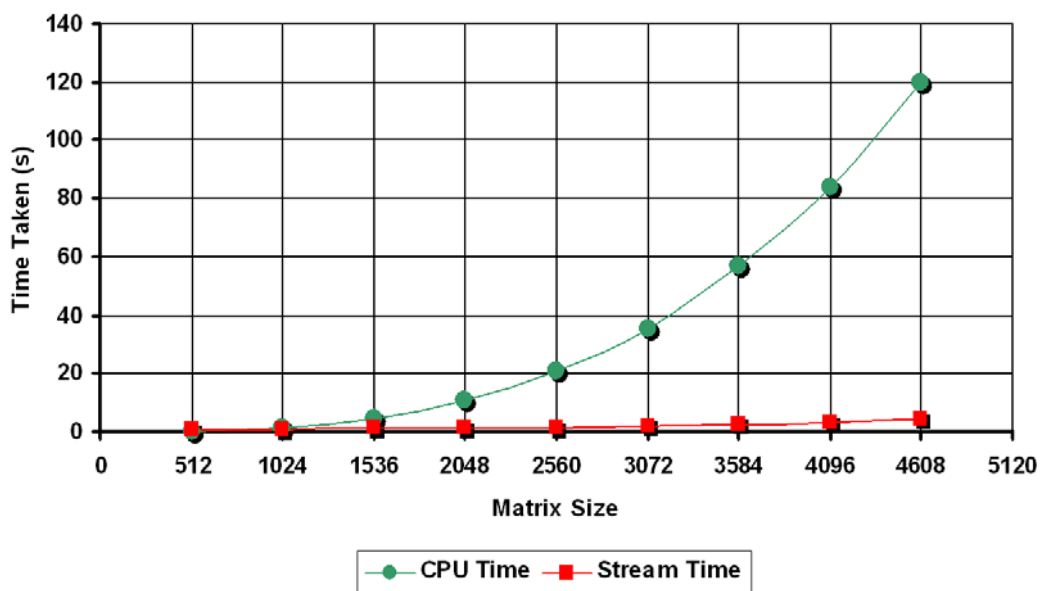


Figure 3 Performance of Brook+ Cholesky Factorization: Stream Processor vs CPU

Figure 4 shows the speedup achieved over the reference CPU implementation. The Radeon™ 4800 Series can achieve speedups of over 25 for this sample implementation.

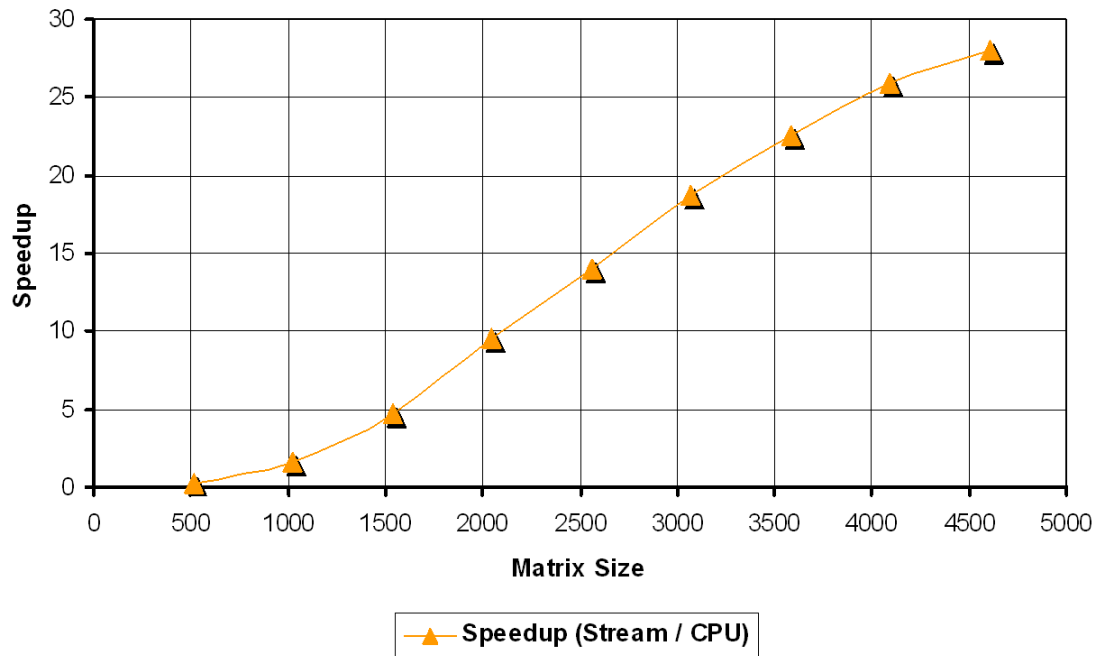


Figure 4 Speedup of Cholesky Factorization: Stream Processor vs CPU Implementation

6 References

1. Jin Hyuk Jung, "Cholesky Decomposition and Linear Programming on a GPU," Scholarly Paper Directed by Dianne P. O'Leary (2006).

Contact

Advanced Micro Devices, Inc.
One AMD Place
P.O. Box 3453
Sunnyvale, CA, 94088-3453
Phone: +1.408.749.4000

For Stream Computing:

URL: <http://ati.amd.com/technology/streamcomputing>

Questions: streamcomputing@amd.com

Developing: streamdeveloper@amd.com

Forum: <http://forums.amd.com/devforum/categories.cfm?catid=38>



The contents of this document are provided in connection with Advanced Micro Devices, Inc. ("AMD") products. AMD makes no representations or warranties with respect to the accuracy or completeness of the contents of this publication and reserves the right to make changes to specifications and product descriptions at any time without notice. The information contained herein may be of a preliminary or advance nature and is subject to change without notice. No license, whether express, implied, arising by estoppel or otherwise, to any intellectual property rights is granted by this publication. Except as set forth in AMD's Standard Terms and Conditions of Sale, AMD assumes no liability whatsoever, and disclaims any express or implied warranty, relating to its products including, but not limited to, the implied warranty of merchantability, fitness for a particular purpose, or infringement of any intellectual property right.

AMD's products are not designed, intended, authorized or warranted for use as components in systems intended for surgical implant into the body, or in other applications intended to support or sustain life, or in any other application in which the failure of AMD's product could create a situation where personal injury, death, or severe property or environmental damage may occur. AMD reserves the right to discontinue or make changes to its products at any time without notice.

Copyright and Trademarks

© 2008 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, ATI, the ATI logo, Radeon, FireStream, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other names are for informational purposes only and may be trademarks of their respective owners.



Printed on
Recycled Paper

Printed in USA